

WHAT IS CLAIMED IS:

1

5

8

10

11

12

13

14

15

16

17

21

24

- A method of identifying one or more positions in a polymer family, the method comprising:
 - (a) accessing data representing a multiple sequence alignment (MSA) of a plurality of polymer sequences; and
- 6 (b) identifying one or more positions within the MSA that have statistically significant conservation energy values using the following equation:

$$\Delta G_i^{stat} = kT^* \sqrt{\sum_{x} \left(\ln \frac{P_i^x}{P_{MSA}^x} \right)^2}$$

9 wherein:

i is a position in the MSA;

 ΔG_i^{stat} is the conservation energy value for position i;

 P_i^x is the probability of monomer x at position i;

 P_{MSA}^{x} is the probability of monomer x in the MSA; and

kT* is an energy unit, where k is Boltzmann's constant.

- 2. The method of claim 1, wherein the method is executed using a machine.
- A program storage device readable by the machine of claim 2 and encoding instructions executable by the machine for performing the operations recited in the claim.
- The method of claim 1, further comprising generating a graphical image of the conservation energy values.
- The method of claim 1, wherein the polymer sequences comprise protein sequences.
- The method of claim 1, wherein monomer x comprises amino acid x.

29

27

Sel	AIT	1
501		2
		3
		4
		5
		6
		7

9

14

15

16

17

18

19

20

21

22

23

- 7. The method of claim 1, wherein the data accessed comprises data from the PDZ domain family.
- 8. The method of claim 1, wherein the data accessed comprises data from the p21^{ras} domain family.
- 7 9. The method of claim 1, wherein the data accessed comprises data from the hemoglobin domain family.
- 10. A method of identifying one or more positions in a polymer family, the method comprising:
- 12 (a) accessing data representing a multiple sequence alignment (MSA) of a
 plurality of polymer sequences;
 - (b) calculating a conservation energy value for each position in the MSA using the following equation:

$$\Delta G_i^{stat} = kT^* \sqrt{\sum_{x} \left(\ln \frac{P_i^x}{P_{MSA}^x} \right)^2}$$

wherein:

i is a position in the MSA;

 ΔG_i^{stat} is the conservation energy value for position i;

 P_i^x is the probability of monomer x at position i;

 P_{MSA}^{x} is the probability of monomer x in the MSA;

kT* is an energy unit, where k is Boltzmann's constant; and

(c) identifying one or more positions within the MSA that have statistically significant conservation energy values.

24 25

11. The method of claim 10, wherein the method is executed using a machine.

27

26

4

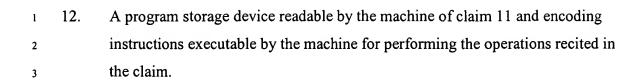
7

10

15

18

21



- 5 13. The method of claim 10, further comprising generating a graphical image of the conservation energy values.
- The method of claim 10, wherein the polymer sequences comprise protein sequences.
- 15. The method of claim 10, wherein monomer x comprises amino acid x.
 - 16. The method of claim 10, wherein the data accessed comprises data from the PDZ domain family.
- 16 17. The method of claim 10, wherein the data accessed comprises data from the p21^{ras} domain family.
- 19 18. The method of claim 10, wherein the data accessed comprises data from the hemoglobin domain family.
- 22 19. A method useful in identifying interacting monomers in a polymer family, the method comprising:
- 24 (a) accessing data representing a multiple sequence alignment (MSA) of a
 25 plurality of polymer sequences;
- 26 (b) calculating a respective conservation energy value for each position in the
 27 MSA using the following equation:

$$\Delta G_i^{stat} = kT^* \sqrt{\sum_{x} \left(\ln \frac{P_i^x}{P_{MSA}^x} \right)^2}$$

wherein:

1		i is a position in the MSA;
2		ΔG_i^{stat} is the conservation energy value for position i;
3		P_i^x is the probability of monomer x at position i;
4		P_{MSA}^{x} is the probability of monomer x in the MSA;
5		kT* is an energy unit, where k is Boltzmann's constant;
6		(c) perturbing a position in the MSA other than position i;
7		(d) re-calculating the respective conservation energy value for each position
8		in the MSA to yield a perturbed conservation energy value; and
9		(e) identifying positions within the MSA that have statistically significant
10		differences between their respective conservation energy values and their
11		perturbed conservation energy values.
12		
13	20.	The method of claim 19, wherein the perturbing includes:
14		selecting a position j in the MSA; and
15		selecting a subset of the MSA, the subset having one or more monomers at
16		position j in the MSA.
17		
18	21.	The method of claim 20, wherein the re-calculating and identifying include:
19		for each position in the MSA, calculating a vector difference $\Delta\Delta G^{\text{stat}}$ between the
20		conservation energy value of the MSA and a conservation energy value of
21		the subset of the MSA using the following equation:
22		$\Delta \Delta G_{i,j}^{stat} = kT^* \sqrt{\sum_{x} \left(\ln \frac{P_{i \partial j}^x}{P_{MSA \partial j}^x} - \ln \frac{P_i^x}{P_{MSA}^x} \right)^2}$
23		wherein:
24		$\Delta\Delta G_{i,j}^{stat}$ is the vector difference in conservation energy values for
25		position i;
26		$P_{i \mathscr{T}}^{x}$ is the probability of monomer x at position i of the subset;
27		$P_{MSA \delta j}^{x}$ is the probability of monomer x in the subset; and

1		identifying positions within the MSA that have statistically significant $\Delta\Delta G^{\text{stat}}$
2		values.
3		
4	22.	The method of claim 21, further comprising generating a graphical image of the
5		$\Delta\Delta G^{\text{stat}}$ values.
6		
7	23.	The method of claim 19, wherein the method is executed using a machine.
8		
9	24.	A program storage device readable by the machine of claim 23 and encoding
10		instructions executable by the machine for performing the operations recited in
11		the claim.
12		
13	25.	The method of claim 19, wherein the polymer sequences comprise protein
14		sequences.
15		
16	26.	The method of claim 19, wherein monomer x comprises amino acid x.
17		
18	27.	The method of claim 19, wherein the data accessed comprises data from the PDZ
19		domain family.
20		
21	28.	The method of claim 19, wherein the data accessed comprises data from the p21 ^{ras}
22		domain family.
23		
24	29.	The method of claim 19, wherein the data accessed comprises data from the
25		hemoglobin domain family.
26		
27	30.	A machine-executed method of quantitatively identifying interacting amino acids
28		in a protein family, the method comprising:
29		(a) accessing data representing a multiple sequence alignment (MSA) of a
30		plurality of protein sequences that are members of a common structural
31		family;

(b) for each position in the MSA, calculating a respective conservation energy value using the following equation:

$$\Delta G_i^{stat} = kT^* \sqrt{\sum_x \left(\ln \frac{P_i^x}{P_{MSA}^x} \right)^2}$$

4 wherein:

i is a position in the MSA;

 ΔG_i^{stat} is the conservation energy value for position i;

 P_i^x is the probability of amino acid x at position i;

 P_{MSA}^{x} is the probability of amino acid x in the MSA;

kT* is an energy unit, where k is Boltzmann's constant;

- (c) selecting a position j in the MSA;
- (d) selecting a subset of the MSA, wherein the subset has one or more amino acids at position j in the multiple sequence alignment;
- (e) for each position in the multiple sequence alignment, calculating a vector difference between the respective conservation energy value of the multiple sequence alignment and the respective conservation energy value of the subset of the multiple sequence alignment; and
- (f) identifying positions within the MSA that have statistically significant vector differences.

19

6

7

8

9

10

11

12

13

14

15

16

17

18

- 20 31. A method of analyzing data comprising:
- 21 (a) providing at least one protein having a crystal structure and multiple
 22 positions;
 - (b) solving the crystal structure of the at least one protein; and
- 24 (c) identifying pathways between interacting positions on the at least one protein.

26

23

- 1 32. A method of analyzing the effect of perturbation on a protein, comprising:
- 2 (a) accessing data representing at least one protein and at least one perturbed 3 protein, both proteins having at least one identical atom;
- 4 (b) calculating a quantity of change Δ_{struct} to the atom using the following equation:

$$\Delta_{struct} = \frac{\left| \vec{r}_{mut} \right|}{\sqrt{\sigma_{mut}^2 + \sigma_{wt}^2}}$$

7 wherein:

6

8

9

10

11

12

13

14

15

16

17

18

19

20

21

 $|\vec{r}_{mut}|$ is the magnitude of a vector connecting the position of the atom in the at least one perturbed protein and the position of the atom in the at least one protein;

 σ_{mut} is a standard deviation of the atom in the at least one perturbed protein; and

 σ_{wt} is a standard deviation of the atom in the at least one protein.

- 33. A method of analyzing data, comprising:
 - (a) accessing data representing at least one protein, a first perturbation of the at least one protein yielding a first perturbed protein, a second perturbation of the at least one protein yielding a second perturbed protein, and a double perturbation of the at least one protein yielding a double perturbed protein, the double perturbation comprising both the first and second perturbations, the proteins each having at least one identical atom;
- 22 (b) calculating a quantity of structural coupling $\Delta\Delta_{struct}$ between the first and second perturbations using the following equation:

$$\Delta\Delta_{struct} = \frac{\left|\vec{r}_{mut1} - \vec{r}_{mut1|mut2}\right|}{\sqrt{\sigma_{wt}^2 + \sigma_{mut1}^2 + \sigma_{mut2}^2 + \sigma_{mut1,mut2}^2}}$$

wherein:

1			\vec{r}_{mut1} is a vector connecting the position of the atom in the first
2			perturbed protein and the position of the atom in the at least
3			one protein;
4			$\vec{r}_{mut1 mut2}$ is a vector connecting the position of the atom in the
5			double perturbed protein and the position of the atom in the
6			second perturbed protein;
7			σ_{wt} is a standard deviation of the atom in the at least one protein;
8			σ_{mut1} is a standard deviation of the atom in the first perturbed
9			protein;
10			σ_{mut2} is a standard deviation of the atom in the second perturbed
11			protein; and
12			$\sigma_{mut1,mut2}$ is a standard deviation of the atom in the double
13			perturbed protein.
14			
15	34.	A me	thod of analyzing microarray data comprising:
16		(a)	accessing microarray data representing an expression level of at least one
17			gene, an expression level of the at least one gene resulting from a first
18			perturbation, an expression level of the at least one gene resulting from a
19			second perturbation, and an expression level of the at least one gene
20			resulting from a double perturbation comprising both the first and second
21			perturbations; and
22		(b)	calculating a degree of coupling $\Delta\Delta E$ between the first and second
23			perturbations using the following equation:
24			$\Delta \Delta E = kT \ln \left(\frac{f_1}{f_2} \right)$
25			wherein:
26			f_1 is the fold effect of the gene due to the first perturbation relative
27			to the at least one gene;

1	f_2 is the fold effect of the gene due to the double perturbation
2	relative to the second perturbation; and
3	kT is an energy unit, where k is Boltzmann's constant.
4	
5	
6	
7	
8	